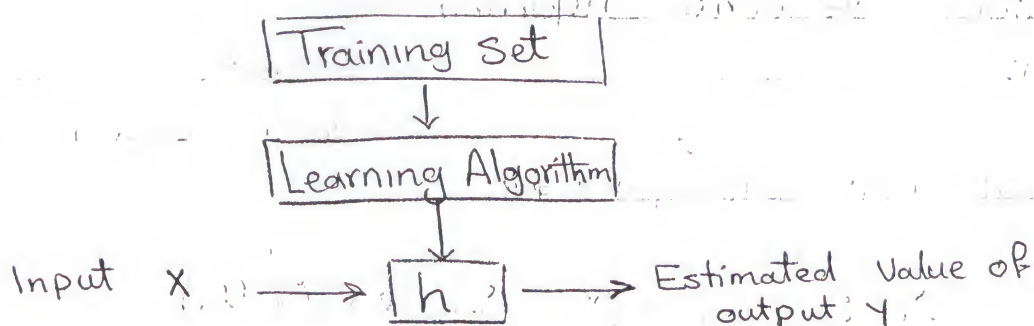


Lec : 4



* Learning Algorithms is used to get Hypothesis

$$h_0(x) = \theta_0 + \theta_1 x$$

* The gradient decent algorithm is Learning Algorithm. used to minimize the Objective Function $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

How

1. Start with some initials θ_0, θ_1
2. Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hope fully end up at a minimum.

example in slide 6:

we start From a θ_0, θ_1 and from these θ_0, θ_1 which we choose to start with, we keep changing them to reduce $J(\theta_0, \theta_1)$ until we reach a minimum point.

example in slide 7:

we start From other values of θ_0, θ_1 and follow the same steps until we reach another minimum point for the same figure.

→ each of these minimum points is called local minimum we want to get the best of them (the minimum of them) and called it global minimum which give the minimum value of $J(\theta_0, \theta_1)$

at First we will study the algorithm of gradient descent generally
(GDA For any Function)

Gradient descent algorithm

Note: it is a general way to get min value of any function so we use it to get h_0
 \downarrow means until end up at a local minimum.

repeat until convergence $\{$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

$\}$

\downarrow Learning rate that controls how big step we take down with gradient descent.

big
we reach to local minimum very fast

small
we reach to local minimum very slow.

So we need to choose suitable α .

- We know that $J(\theta_0, \theta_1)$ is constant for each line so when we change θ_0, θ_1 together or change θ_0 only or change θ_1 only we get another different line with new $J(\theta_0, \theta_1)$.
- During gradient descent algorithm we must change θ_0, θ_1 together in each step (at the same time) to get one value to $J(\theta_0, \theta_1)$ of only one line.

Correct:

one move $\left[\begin{array}{l} \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_0 := \text{temp0} \\ \theta_1 := \text{temp1} \end{array} \right]$

$J(\theta_0, \theta_1)$ for one line
 θ_0, θ_1 change in the same time.

Incorrect

2 moves $\left[\begin{array}{l} \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_0 := \text{temp0} \\ \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_1 := \text{temp1} \end{array} \right]$

$J(\theta_0, \theta_1)$ for one line in which θ_0 changes but θ_1 const.

$J(\theta_0, \theta_1)$ for another line in which θ_0 constant and θ_1 changes.

• Simplified ($\phi_0 \equiv 0$)

Hypothesis: $h_{\phi}(x) = \phi_1 x$

Parameters: ϕ_1

Cost Function: $J(\phi_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\phi}(x^{(i)}) - y^{(i)})^2$

Goal: minimize $J(\phi_1)$ (minimize error to get the Best line)
 ϕ_1

Gradient descent algorithm. (one of many ways that can be used to get $h_{\phi}(x)$)
→ we use it to minimize error $J(\phi_0, \phi_1)$ to get the best Line $h_{\phi}(x)$
repeat until Convergence &

$$\phi_1 := \phi_1 - \alpha \frac{d}{d\phi_1} J(\phi_1)$$

}

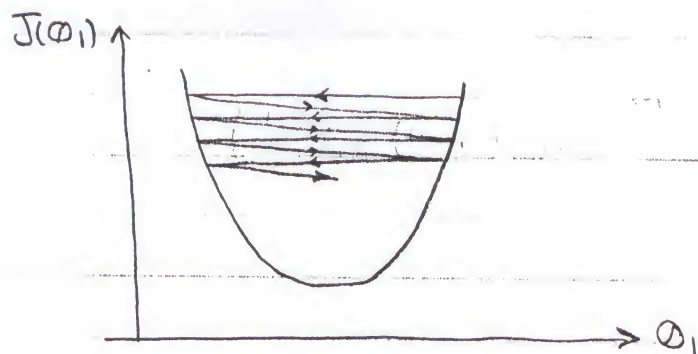
So → Two Parameters Control GDA ($\alpha, \frac{d}{d\phi_1} J(\phi_1)$)
to get minimum $J(\phi_1)$

Alpha α

- IF α is too Small, gradient descent can be slow.
- IF α is too Large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.
- Gradient descent can ^{تقارب} converge to local minimum, even with the learning rate α Fixed.
- As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.

Neglecting the effect of the slope ($\frac{dJ(\theta_1)}{d\theta_1}$)

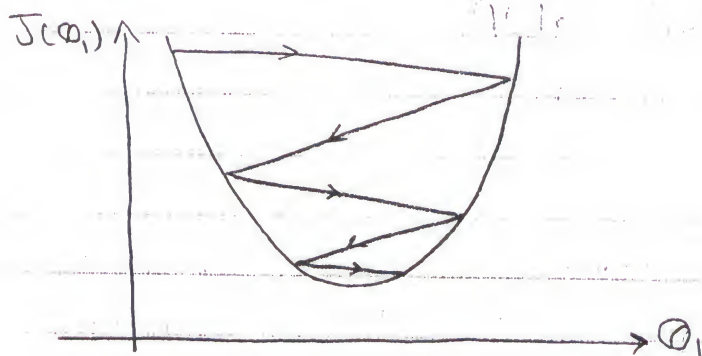
① α is too small



→ The values of θ_1 will move between big and small values (+ve & -ve values of slope)

→ It will reach to Local minimum point but very slow.

② α is too great



→ The values of θ_1 will move between big and small values (+ve & -ve values of slope)

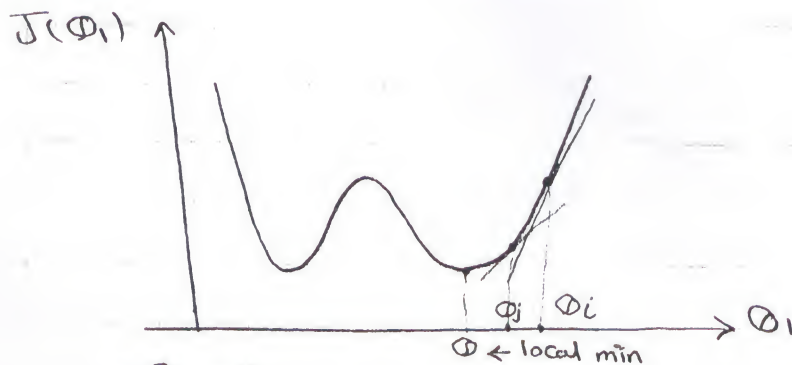
→ It may not be able to reach to Local minimum because of its Large step.

2 $\frac{dJ(\phi_1)}{d\phi_1} \leftarrow \text{Slope}$

(المشتقة الأولى في حساب التفاضل)

* Assuming That α is Fixed

① Slope = +ve Value



$$\phi_{i_{\text{new}}} := \phi_{i_{\text{old}}} - \alpha \frac{dJ(\phi_1)}{d\phi_1}$$

$\rightarrow \alpha$ Fixed

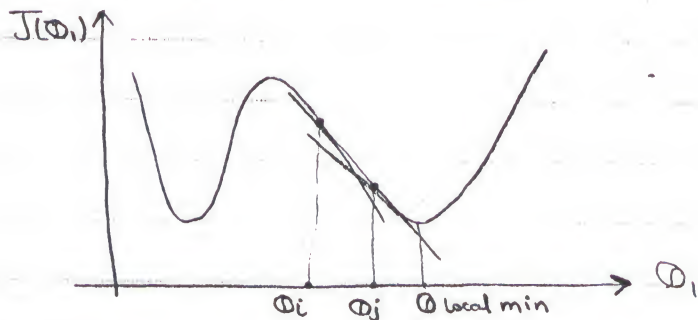
$\rightarrow \phi_{i_{\text{new}}}$ depends on $\frac{dJ(\phi_1)}{d\phi_1}$

* $\phi_{i_{\text{old}}} = \phi_i$

$$\phi_{i_{\text{new}}} = \phi_i - \text{Const} * +ve = \phi_j \rightarrow \phi_{i_{\text{new}}} \text{ will decrease}$$

\rightarrow we will reach to Local minimum.

② Slope = -ve Value.



* $\phi_{i_{\text{old}}} = \phi_i$

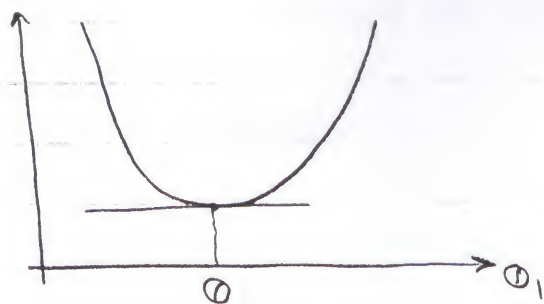
$$\phi_{i_{\text{new}}} = \phi_i - (\text{const} * -ve) = \phi_i + (\text{const} * +ve) = \phi_j$$

$\rightarrow \phi_{i_{\text{new}}} = \phi_j$ will increase

\rightarrow we will reach to Local minimum.

Note : • the value of slope doesn't effect reaching to Local min
it only effects the value 2ϕ of ϕ_1 by increase or decrease.

→ at Local minimum.



$$\theta_{\text{new}} = \theta = \theta_{\text{old}} - \alpha \frac{dJ(\theta_1)}{d\theta_1}$$

at this point $\frac{dJ(\theta_1)}{d\theta_1} = \text{Zero}$

∴ $\theta_{\text{new}} = \theta_{\text{old}}$ ∴ θ_1 doesn't change

→ We reach to the Final Point when $J(\theta_1)$, which depends on the value change of θ_1 , doesn't change.
"Here we reach to Local minimum."

Note: the Slope Change changes the theta θ_1 value (increase, decrease) but doesn't change our way to reach Local minimum.

Gradient Descent For Linear Regression With One Variable.

→ Here we will study how to use gradient descent algorithm to get the Best Linear regression model $h_\theta(x) = \theta_0 + \theta_1 x$ by minimizing the objective Function $J(\theta_0, \theta_1)$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

→ Gradient descent algorithm.

repeat until convergence &

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

} (For $j=1$ & $j=0$)

to use gradient descent algorithm to get best Linear regression we need to get $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ to be able to use this algorithm.

→ Gradient Descent Algorithm For Linear Regression :-

repeat until convergence &

$$\theta_0 := \theta_0 - \alpha \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$ حيث θ_1 ثابتة

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$ حيث θ_0 ثابتة

→ update θ_0 and θ_1 Simultaneously

قيمه θ_0 , θ_1 تتغيران في تزامن كما سبق الشرح

→ Using Matlab to draw GDA For linear regression we get the draw in slide 17

We Note that it has one Local minimum point (global minimum)

So, This mean that IF we start From any value of θ_0, θ_1 we reach to the same point

by another meaning :

IF we start From known values of θ_0, θ_1 For the First time and start From different values of θ_0, θ_1 From th First time For the second time and repeate starting From different values For every time we reach to the same Local minimum point

. Here, we have one Local minimum point we can called it global minimum.

Local minimum \equiv Global minimum ..

"Batch" Gradient Descent

Batch : Each step of gradient descent uses all the training examples.